

## Regular languages and finite automata.

Suppose  $\Sigma$  is a finite alphabet.

Recall the following basic operations on languages

$$L, K \subseteq \Sigma^*$$

### Concatenation

$$LK = \{xy \mid x \in L, y \in K\}$$

### Union

$$L \cup K = \{x \mid x \in L \text{ or } x \in K\}$$

### Kleene star

$$L^* = L^0 \cup L^1 \cup L^2 \cup \dots$$

$$= \{w_1 w_2 \dots w_n \mid n \geq 0 \text{ \& } w_i \in L\}.$$

Definition The class of regular languages in  $\Sigma$  is the smallest class  $\mathcal{R}$  of languages containing

$\emptyset$ ,  $\{a\}$  (for  $a \in \Sigma$ ) and such that if

$L, K \in \mathcal{R}$ , then also  $LK, L \cup K, L^* \in \mathcal{R}$ .

Example All finite languages are regular, since

$$\{a_1 a_2 \dots a_n\} = \{a_1\} \circ \{a_2\} \circ \dots \circ \{a_n\}.$$

Definition A deterministic finite state automaton, DFA, consists of a finite directed graph  $M$ , i.e.,  $M = (V, E)$  where  $V$  is a finite set of vertices and  $E \subseteq V^2$  is a set of directed edges, along with

(1) a distinguished start state  $s_0 \in V$ ,

(2) a set  $A \subseteq V$  of accepting states,

(3) a labelling  $l: E \rightarrow \mathcal{P}(\Sigma) \setminus \{\emptyset\}$

for any  $s \in V$  and any  $a \in \Sigma$  there is exactly one edge  $(s, t) \in E$  originating at  $s$  and with label  $a$ , i.e.,  $a \in l(s, t)$ .

We call  $V$  the set of states of the automaton.

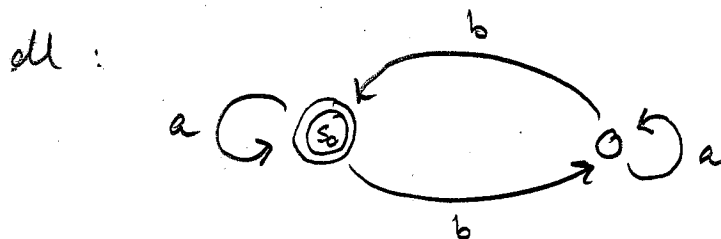
Given an automaton  $M$  as above and a string

$w = a_1 a_2 \dots a_n \in \Sigma^*$ , we say that  $M$  accepts  $w$  if

the unique edgepath  $(e_1, e_2, \dots, e_n)$  in  $M$  originating at  $s_0$  and with  $a_j \in l(e_j)$  terminates at an accepting state.

Given  $M$ , we let  $L(M)$  be the language consisting of the strings accepted by  $M$ . We also say that  $M$  recognizes  $L(M)$ .

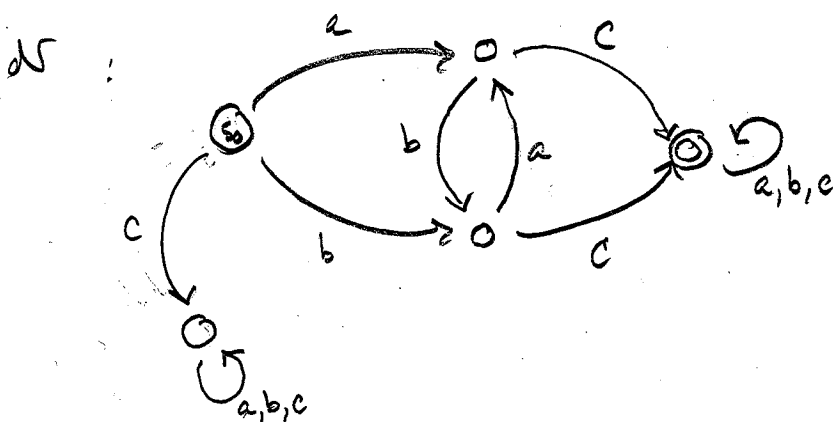
Example We draw a DFA by a diagram in the plane as a usual directed graph with labeled edges. Moreover, the accepting states are indicated by double circles.



Here dl has exactly two states and the only accepting state is  $s_0$ . We see that  $\Sigma = \{a, b\}$  and

$$L(dl) = \{w \mid |w|_b = \# \text{ of occurrences of } b \text{ in } w \text{ is even}\}.$$

Similarly, for  $\Sigma = \{a, b, c\}$



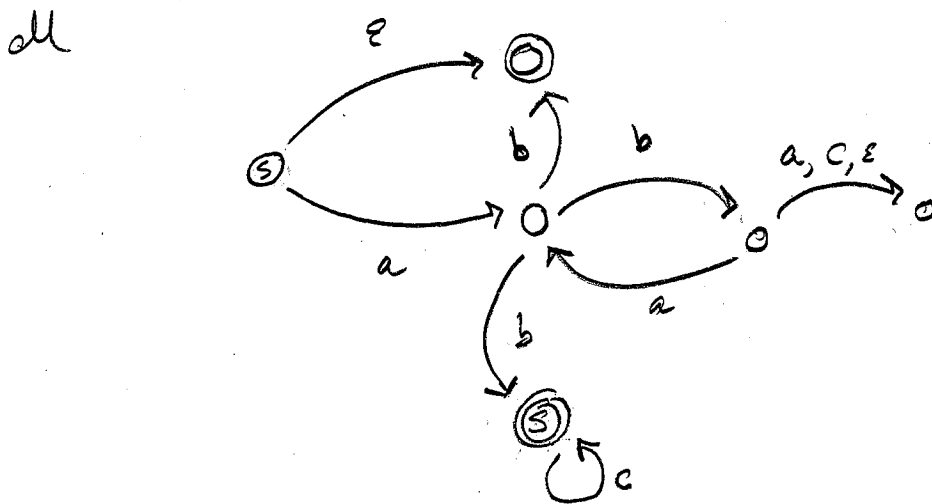
$$L(dV) = \{awcv, bwcv \mid w \in \{a, b\}^*, v \in \{a, b, c\}^*\}.$$

Definition A generalised non-deterministic finite state automaton,  $\epsilon$ -NFA, is a directed finite graph  $M = (V, E)$  along with

- (1) a non-empty set of start states  $S \subseteq V$ ,
- (2) a set of accepting states  $A \subseteq V$ ,
- (3) a labelling  $l: E \rightarrow \mathcal{P}(\Sigma \cup \{\epsilon\}) \setminus \{\emptyset\}$ .

Given an  $\epsilon$ -NFA  $M$  and a string  $w \in \Sigma^*$ , we say that  $M$  accepts  $w$  if there is an edge path  $(e_1, \dots, e_m)$  and labels  $b_i \in l(e_i)$  ( $b_i$  can be  $\epsilon$ ) such that  $w = b_1 b_2 \dots b_m$ .

Example We indicate the start states by an  $S$ .



$$L(M) = \{ (ab)^n c^m \mid n \geq 0, m \geq 0 \}$$

## Conc types

Suppose  $L \subseteq \Sigma^*$  is a language and  $w \in \Sigma^*$ .

We define the conc type of  $w$  wrt  $L$  by

$$\text{conc}_L(w) = \{x \in \Sigma^* \mid wx \in L\}$$

and let the conc types of  $L$  be

$$\text{Conc}(L) = \{\text{conc}_L(w) \mid w \in \Sigma^*\}$$

Note The sets of conc types of  $L$  are the concs of all words  $w \in \Sigma^*$ , not just  $w \in L$ .

Example Let  $L = \{a^n b^n \mid n \geq 0\} \subseteq \{a, b\}^*$

$$\text{Then } \text{conc}_L(b^2) = \emptyset, \quad \text{conc}_L(\varepsilon) = L,$$

$$\text{conc}_L(a^2) = \{a^n b^{n+2} \mid n \geq 0\}.$$

Note also that  $\text{conc}_L(a), \text{conc}_L(a^2), \dots$  are all distinct.

Definition For  $L \subseteq \Sigma^*$  a language define an equivalence relation  $\sim_L$  on  $\Sigma^*$  by

$$w \sim_L v \iff w \text{ and } v \text{ have the same conc type, i.e., } \text{conc}_L(w) = \text{conc}_L(v).$$

Lemma Let  $M$  be an  $\epsilon$ -NFA and let  $L = L(M)$  be the language accepted by  $M$ . Then  $L$  has only finitely many cone types, i.e.,  $\nu_2$  has only finitely many classes.

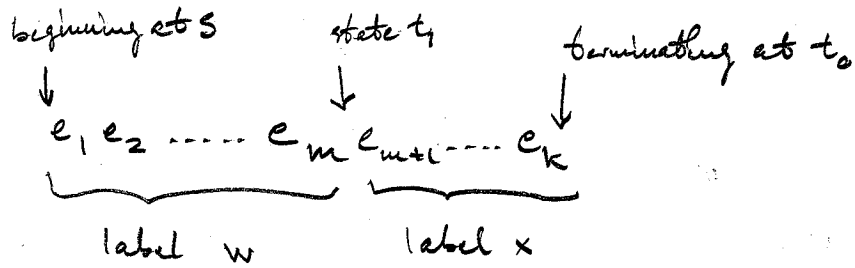
Proof We define another equivalence relation  $\approx$  on  $\Sigma^*$  as follows:

$w \approx v \iff$  for any start state  $s$  and arbitrary state  $t$  of  $M$ , there is an edge path from  $s$  to  $t$  with edge label  $w$  if and only if there is an edge path from  $s$  to  $t$  with edge label  $v$ .

Note then that if  $M$  has  $n$  states, then  $\approx$  has at most  $2^{n^2}$  classes (for each edge  $(s, t)$ , we have to respond to a yes-no question).

Claim If  $w \approx v$  then also  $w \nu_2 v$ .

For suppose that, e.g.,  $x \in \text{cone}_2(w)$ . Then there is an edge path beginning at a start state  $s$  and terminating at an accepting state  $t_0$  with edge label  $wx$ . Let  $t_1$  be any state along this edge path such

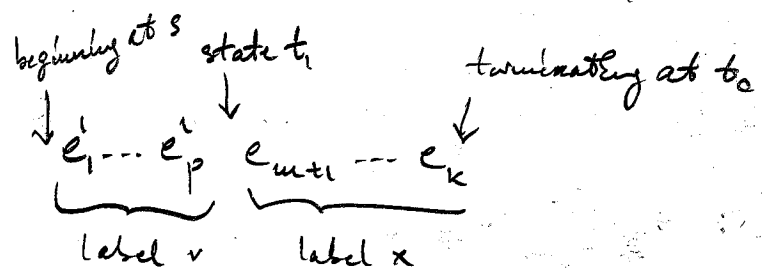


that the edge path arrives at  $t_1$  with label  $w$ .

Then there is an alternative edge path

$e'_1 \dots e'_p$  beginning at  $s$  and terminating at  $t_1$  with label  $v$ . It follows that there is an

edge path from  $s$  to  $t_0$  with label  $vx$ , where  $x \in cone_L(v)$ .



Similarly,  $cone_L(v) \subseteq cone_L(w)$ . □

Lemma Suppose  $L \subseteq \Sigma^*$  is a language with finitely many cone types. Then there is a DFA  $M$  with  $L(M) = L$ .

Proof Let  $\Delta_1, \dots, \Delta_n$  be the finitely many cone types of  $L$  and note that for  $w, x \in \Sigma^*$  and  $a \in \Sigma$ , if  $w \sim_L x$  then also  $wa \sim_L xa$ .

Let  $M$  have states  $\Delta_1, \dots, \Delta_n$  and put an arrow  $\Delta_i \xrightarrow{a} \Delta_j$  if for some

any  $w \in \Sigma^*$  with  $\text{conv}_2(w) = \Delta_i$ , we have

$\text{conv}_2(wa) = \Delta_j$ . Then a state  $\Delta_i$  is accepting

if  $\varepsilon \in \Delta_i$  and  $s = \text{conv}_2(\varepsilon)$  is the unique

start state.  $M$  is clearly deterministic and

$L(M) = L$ .  $\square$

Lemma Let  $M$  be a DFA. Then the language

$L(M)$  is regular.

Proof Let  $V$  be the finite set of states of  $M$  and let  $t_0, t_1 \in V$  be arbitrary. For any

$X \subseteq V$ , let

$G(X, t_0, t_1) = \{w \in \Sigma^* \mid \text{there is an edge path from } t_0 \text{ to } t_1 \text{ only passing through states in } X \text{ and having label } w\}$ .

By induction on  $|X|$ , we show that for any  $t_0, t_1$ , the language  $G(X, t_0, t_1)$  is regular.

$|X| = 0$ : In this case  $X = \emptyset$  and so  $w \in G(X, t_0, t_1)$

if and only if  $w = a$  for some  $a \in \Sigma$  for which there is an edge  $t_0 \xrightarrow{a} t_1$ .

So  $G(X, t_0, t_1)$  is a (finite) subset of  $\Sigma$

and hence is regular.

$|X| = n+1$  : Assume the result holds for all subsets of  $V$  of size  $\leq n$  and assume  $|X| = n+1$ .

Then we have that

$$G(X, t_0, t_1) = \left( \bigcup_{q \in X} G(X, \{q\}, t_0, t_1) \right)$$

$$\cup \left( \bigcup_{q \in X} G(X, \{q\}, t_0, q) \circ G(X, \{q\}, q, q)^* \circ G(X, \{q\}, q, t_1) \right)$$

which is regular by the induction hypothesis.

Clearly, the right hand side is contained in  $G(X, t_0, t_1)$ .

Conversely, suppose  $w \in G(X, t_0, t_1)$  and consider the edge path from  $t_0$  to  $t_1$  with label

$$w = a_1 a_2 \dots a_n$$

$$t_0 \xrightarrow{a_1} q_1 \xrightarrow{a_2} q_2 \longrightarrow \dots \xrightarrow{a_{n-1}} q_{n-1} \xrightarrow{a_n} t_1$$

Then, if  $|w| = n \leq 1$ , we have  $w \in \bigcup_{q \in X} G(X, \{q\}, t_0, t_1)$ .

Otherwise, note that

$$a_1 \in G(X, \{q_1\}, t_0, q_1)$$

• if  $q_i$  is the last occurrence of  $q_i$  among  $q_1, \dots, q_{n-1}$ , then

$$a_{i+1} \dots a_n \in G(X, \{q_i\}, q_i, t_i)$$

- if  $q_i = q_j = q_1$  for  $i < j$  and  $q_l \neq q_1$   
for all  $i < l < j$ , then

$$a_i a_{i+1} \dots a_{j-1} \in G(X, \{q_1\}, q_1, q_1)$$

Thus,  $w = a_1 \dots a_n$  belongs to

$$G(X, \{q_1\}, t_0, q_1) \circ G(X, \{q_1\}, q_1, q_1)^* \circ G(X, \{q_1\}, q_1, t_1) \quad \blacktriangle$$

So also  $L(M) = \bigcup_{q \in A} G(N, s_0, q)$  is regular. □

Theorem TFAE for a language  $L$  over a finite alphabet  $\Sigma$ .

(a)  $L$  is regular

(b)  $L = L(M)$  for some DFA  $M$ .

(c)  $L = L(M)$  for some  $\epsilon$ -NFA  $M$

(d)  $\text{Comp}(L)$  is finite.

There the equivalence of regular languages and languages recognized by finite automata is known as the Kleene theorem, while the equivalence with (d) is the Myhill-Nerode theorem.

Proof We have already proved  $(c) \Rightarrow (d) \Rightarrow (b) \Rightarrow (a)$ .

So we need only prove that regular languages are of the form  $L(M)$  for  $\epsilon$ -NFA  $M$ .

Since one can easily build  $\epsilon$ -NFA recognizing any finite language, it suffices to show that if  $L, K$  are recognized by DFA, then also  $LK$ ,  $L \cup K$  and  $L^*$  are recognized by  $\epsilon$ -NFA.

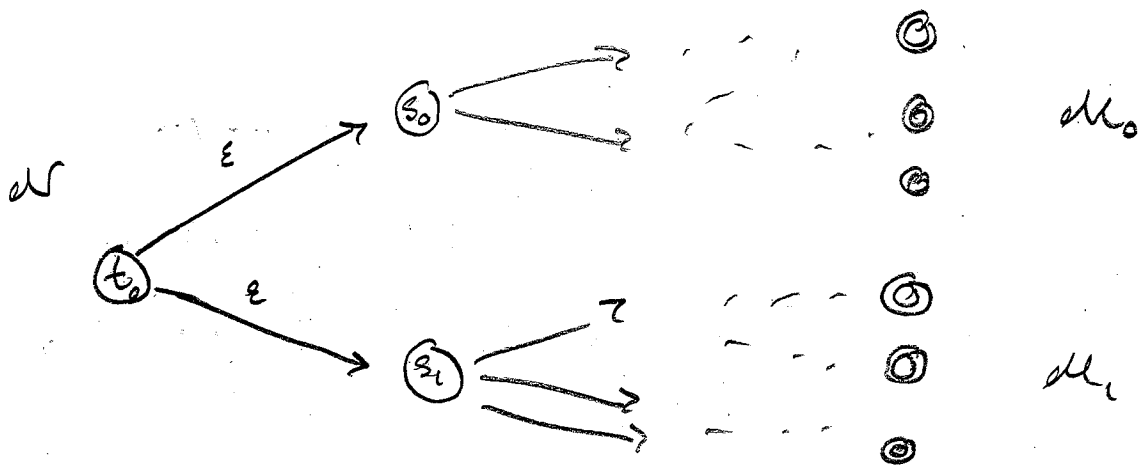
So suppose  $M_0 = (V_0, E_0, s_0, A_0, l_0)$  and  $M_1 = (V_1, E_1, s_1, A_1, l_1)$  are deterministic finite automata recognizing  $L$  and  $K$  respectively.

That is,  $V_i$  are the states,  $E_i \subseteq V_i^2$  are directed edges,  $s_i$  are start states,  $A_i$  are accepting states and  $l_i: E_i \rightarrow P(\Sigma) \setminus \{\emptyset\}$  are labellings. Wlog,  $V_0 \cap V_1 = \emptyset$ .

Let us first build an  $\epsilon$ -NFA recognizing  $L \cup K$ :

Let  $M = (V_0 \cup V_1 \cup \{t_0\}, E_0 \cup E_1 \cup \{(t_0, s_0), (t_0, s_1)\}, t_0, A_0 \cup A_1, l)$

where  $l(e) = l_0(e)$  for  $e \in E_0$ ,  $l(e) = l_1(e)$  for  $e \in E_1$  and  $l(t_0, s_0) = l(t_0, s_1) = \epsilon$ .



Thus, at the first stage of the computation, w/o reading any input of the string,  $N$  has to decide to feed the input to either  $M_0$  or  $M_1$ . So  $L(N) = L \cup K$ .

Now, let us construct  $N$  to recognize  $LK$ .

$$N = (V_0 \cup V_1, E_0 \cup E_1 \cup \{(q, s_1) \mid q \in A_0\}, s_0, A_1, \ell)$$

where  $\ell(e) = \ell_i(e)$  for  $e \in E_i$  and  $\ell(q, s_1) = \epsilon$  for  $q \in A_0$ . Then, on input  $w$ ,  $N$  begins with a computation in  $M_0$  and can from any accepting state of  $M_0$  jump to  $s_1$  and then continue with a computation in  $M_1$ . So  $N$  accepts  $w$  if and only if  $w = xy$ , where  $x \in L$  and  $y \in K$ . Thus,  $L(N) = LK$ .

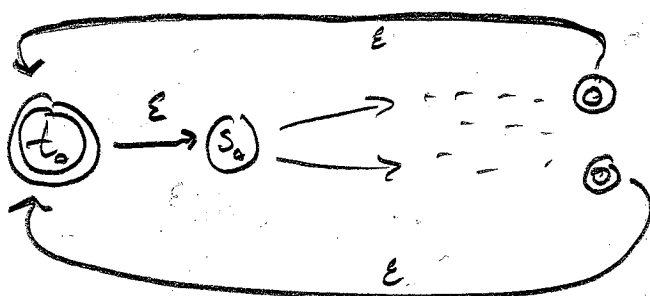
Finally, for  $L^*$  let  $t_0$  be a new state and set

$$\mathcal{N} = (V_0 \cup \{t_0\}, E_0 \cup \{(t_0, \varepsilon)\} \cup \{(q, t_0) \mid q \in A_0\}, t_0, A_0 \cup \{t_0\}, \delta)$$

where  $\delta(e) = \delta_0(e)$  for  $e \in E_0$ ,  $\delta(t_0, \varepsilon) = \varepsilon$ ,

$\delta(q, t_0) = \varepsilon$  for  $q \in A_0$ .

Again,  $L(\mathcal{N}) = L^*$ .



□

Remark We can of course also show that regular languages are recognized by DFA by instead using Myhill-Nerode theorem, i.e., by showing that they only have finitely many cone types.

For example, suppose  $L$  and  $K$  have finitely many cone types each and let  $\Delta_1, \dots, \Delta_n$  be the cone types of  $K$ .

Now set

$$w \approx v \iff w \sim_L v$$

and for any  $i \leq n$ ,  $w$  has a decomposition as  $w = xy$  with  $x \in L$  and  $\text{cone}_K(y) = \Delta_i$  if and only if  $v$  has.

Clearly,  $\approx$  is an equivalence relation with finitely many classes (note that the second part is induced by  $n$  yes/no questions). Moreover, suppose

that  $w \approx v$  and that  $wu \in LK$ . Then

either we can write  $u = xy$ , where  $y \in K$  and  $wx \in L$ , whence also  $vx \in L$  and thus

$wu = vxy \in LK$ , or we can write  $w = xy$ , where  $x \in L$  and  $yu \in K$ . In the second case,

let  $\Delta_i = \text{conc}_K(y)$  and note that then

$v$  has a decomposition as  $v = st$ , where

$s \in L$  and  $\text{conc}_K(t) = \Delta_i$ . It thus follows that

$tu \in K$ , whence  $vu = sttu \in LK$ .

In any case,  $vu \in LK$ , so  $\text{conc}_{LK}(w) \subseteq \text{conc}_{LK}(v)$ .

By symmetry,  $w \approx_{LK} v$ . So  $\approx$  reduces to  $\approx_{LK}$  and hence  $LK$  has finitely many conc types.

Exercise: Show that if  $L, K$  have finitely many conc types, then also  $L \cup K$  and  $L^*$  have finitely many conc types.

Before giving more examples of regular languages, we give some indication of their limitations.

## Pumping lemma

Let  $L \subseteq \Sigma^*$  be a regular language. Then there is an integer  $n \geq 1$  such that any word  $w \in L$  with  $|w| > n$  can be expressed as

$$w = xyz, \quad y \in \Sigma^+, \quad x, z \in \Sigma^*$$

where  $|x| \leq n$  and  $xy^i z \in L$  for all  $i \geq 0$ .

### Proof

First suppose  $L$  is recognized by a DFA  $M$  with  $n$  states. Then if  $w$  is a word of length  $> n$  the edge path labeled by  $w$  will have to pass through some state twice. So let  $x$  be the shortest prefix of  $w$  at which a state  $q$  is reached that is being visited twice. Clearly  $|x| \leq n$ . Also, let  $y$  be the shortest non-empty string st. the edge path with label  $xy$  arrives to  $q$ . Then so does  $xy^i$  for all  $i \geq 0$  and thus  $xy^i z \in L$  for all  $i \geq 0$ .  $\square$

Example The languages  $\{a^n b^n \mid n \geq 0\}$  and  $\{a^n \mid n \text{ is a prime}\}$  are not regular.

## The subsequence ordering

Definition Let  $<$  be a strict partial ordering on a set  $X$ , i.e.,  $<$  is transitive and irreflexive. We say that  $<$  is a well-quasi-ordering, wqo, if

(i) any antichain is finite, i.e., in any infinite sequence  $(x_i)_{i \in \mathbb{N}}$  there are  $i \neq j$  st.  $x_i = x_j$  or  $x_i < x_j$ ,

(ii) there is no infinite descending chain, i.e., no infinite sequence  $(x_i)_{i \in \mathbb{N}}$  with  $x_1 > x_2 > x_3 > \dots$  ( $<$  is well-founded)

Exercise Set  $x \leq y \Leftrightarrow x < y$  or  $x = y$ .

Show that  $<$  is wqo if and only if

for any infinite sequence  $(x_i)_{i \in \mathbb{N}}$  there are  $i < j$  with  $x_i \leq x_j$ .

[Hint: Use Ramsey's theorem].

Definition Suppose  $w = a_1 a_2 \dots a_n$  and  $x$  are words in  $\Sigma$ . We say that  $w$  is a subsequence of  $x$  if there are words  $y_0, \dots, y_n \in \Sigma^*$  with  $x = y_0 a_1 y_1 a_2 y_2 \dots a_n y_n$  and  $w \neq x$ . Write  $w < x$  to denote this.

Proposition Let  $\Sigma$  be a finite alphabet.

Then  $(\Sigma, <)$  is a well-ordered set.

Proof Assume towards a contradiction that there is some sequence  $(x_i)_{i \in \mathbb{N}}$  st.  $i < j \Rightarrow x_i \not\leq x_j$  and call any such sequence bad.

We inductively construct a bad sequence as follows:

• Let  $y_1 \in \Sigma^*$  be a word of shortest length beginning an infinite bad sequence.

• Let  $y_2 \in \Sigma^*$  be a word of shortest length st.  $y_1, y_2$  begins an infinite bad sequence.

• Let  $y_3 \in \Sigma^*$  be a word of shortest length beginning an infinite bad sequence, etc.

Then, as  $\Sigma$  is finite, there is an infinite subsequence, say  $(y_{n_i})_{i \in \mathbb{N}}$ , with constant first letter, eq.  $y_{n_i} = a z_{n_i}$ . Note again that  $z_{n_i} \not\leq z_{n_j}$  for  $i < j$  and  $z_{n_i} < y_{n_i}$ , so

$y_1, y_2, \dots, y_{n_1-1}, z_{n_1}, z_{n_2}, \dots$

is also an infinite bad sequence, but

with  $|z_{n_1}| < |y_{n_1}|$ , contradicting the minimality of

$y_1, y_2, \dots, y_{n_1-1}, y_{n_1}, y_{n_1+1}, \dots$

□

Definition Let  $(X, <)$  be a strict partial ordering and  $B \subseteq Y \subseteq X$  subsets. We say that  $B$  is a basis for  $Y$  if

$$Y = \{x \in X \mid \exists z \in B \ z \leq x\}.$$

Exercise Show that  $(X, <)$  is wqo if and only if any  $Y \subseteq X$ , which is closed upwards, i.e.,  $(y \leq x \ \& \ y \in Y) \Rightarrow x \in Y$ , has a finite basis.

Corollary Let  $L \subseteq \Sigma^*$  be a language closed under supersequences, i.e., if  $x \in L$  and  $x \prec y$ , then  $y \in L$ .

Then  $L$  has a finite basis  $B = \{x_1, \dots, x_n\}$ ,

$$\text{i.e., } L = \{y \in \Sigma^* \mid x_i \leq y \text{ for some } i=1, \dots, n\}.$$

Theorem Let  $L \subseteq \Sigma^*$  be any language closed under supersequences. Then  $L$  is regular.

Proof Let  $B = \{x_1, \dots, x_n\}$  be a finite basis for  $L$  and let  $C$  be the finite set of all prefixes of elements in  $B$ .

Set  $w \approx v \Leftrightarrow$  for any  $x \in C$ ,  $x \leq w$  if and only if  $x \leq v$ .

Then  $w \approx v \Rightarrow w \sim_2 v$ , and since  $\approx$  has only finitely many classes, so does  $\sim_2$ , whence  $L$  is regular.  $\square$

Theorem Let  $L \subseteq \Sigma^*$  be regular, then so is  $\Sigma^* \setminus L$ .

PF Note that  $\text{cove}_{\Sigma^* \setminus L}(w) = \Sigma^* \setminus \text{cove}_L(w)$ , so if  $L$  has only finitely many cove types, the same holds for  $\Sigma^* \setminus L$ .  $\square$

Corollary Let  $L \subseteq \Sigma^*$  be a language closed under taking subsequences. Then  $L$  is regular.

PF Just note that  $\Sigma^* \setminus L$  is closed under supersequences, so  $\Sigma^* \setminus L$  and thus also  $L$  are regular.  $\square$

Example Let  $\Sigma = \{0, 1, \dots, 9\}$  and let  $L \subseteq \Sigma^+$  be the set of all prime numbers written in base 10 and  $K = \{x \in \Sigma^+ \mid \exists y \in L \ y \leq x\}$ . Then  $K$  has a finite basis  $B = \{x_1, \dots, x_n\}$ , where  $x_1, \dots, x_n$  are actually prime numbers written in base 10.

Example Let  $\Sigma = \{0, 1, 2\}$  and let  $L \subseteq \Sigma^+$

be the set of all prime numbers written in base 3. Then  $\{2, 10, 111\}$  is a basis for

$$\text{pref}(L) = \{y \in \Sigma^+ \mid \exists x \in L \ x \preceq y\}.$$

To see this, note first that  $2 \sim 2$ ,  $10 \sim 3$ ,  $111 \sim 13$  (where 2, 3, 13 are in base 10). So  $2, 10, 111 \in L$ .

Also, suppose  $x \in L$ . Then if  $2 \not\preceq x$ , we have

$x \in \{0, 1\}^*$ , and if further  $10 \not\preceq x$ , also

$x \in 0^*1^*$ . Now, unless  $x$  represents the number

0,  $x$  has no leading 0 and hence  $x \in 1^*$ .

So suppose  $x \in 1^*$ . Then as neither 1, nor 11 represent prime numbers, we must have

$x \in 1111^*$ , i.e.,  $111 \preceq x$ .

Notation Let  $x \equiv y$  if  $x$  is a prefix of  $y$ .

Theorem Suppose  $L$  is a regular language.

Then  $\text{pref}(L) = \{x \in \Sigma^* \mid \exists y \in L \ x \equiv y\}$

is regular too.

Proof Exercise. □